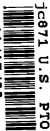


11/15/00



11/15/00

11-16-00

A

Please type a plus sign (+) inside this box → ☒

Approved for use through 09/30/2000. OMB 0651-0032
 Patent and Trademark Office: U.S. DEPARTMENT OF COMMERCE
 Under the Paperwork Reduction Act of 1995, no persons are required to respond to a collection of information unless it displays a valid OMB control number.

UTILITY PATENT APPLICATION TRANSMITTAL

(Only for new nonprovisional applications under 37 C.F.R. § 1.53(b))

Attorney Docket No. 2942-991842

First Inventor or Application Identifier Mark Kantrowitz

Title Method And Apparatus For Efficient Identification Of Duplicate ...

Express Mail Label No. EL561512695US

APPLICATION ELEMENTS

See MPEP chapter 600 concerning utility patent application contents.

1. ☒ * Fee Transmittal Form (e.g., PTO/SB/17)
 (Submit an original and a duplicate for fee processing)
2. ☒ Specification [Total Pages 19]
 (preferred arrangement set forth below)
 - Descriptive title of the Invention
 - Cross References to Related Applications
 - Statement Regarding Fed sponsored R & D
 - Reference to Microfiche Appendix
 - Background of the Invention
 - Brief Summary of the Invention
 - Brief Description of the Drawings (if filed)
 - Detailed Description
 - Claim(s)
 - Abstract of the Disclosure
3. ☒ Drawing(s) (35 U.S.C. 113) [Total Sheets 9]
 4. Oath or Declaration [Total Pages 2]
 a. ☒ Newly executed (original or copy)
 b. ☐ Copy from a prior application (37 C.F.R. § 1.63(d))
 (for continuation/divisional with Box 16 completed)
 i. ☐ DELETION OF INVENTOR(S)
 Signed statement attached deleting inventor(s) named in the prior application, see 37 C.F.R. §§ 1.63(d)(2) and 1.33(b).

* NOTE FOR ITEMS 1 & 3: IN ORDER TO BE ENTITLED TO PAY SMALL ENTITY FEES, A SMALL ENTITY STATEMENT IS REQUIRED (37 C.F.R. § 1.27), EXCEPT IF ONE FILED IN A PRIOR APPLICATION IS RELIED UPON (37 C.F.R. § 1.28).

ADDRESS TO:

 Assistant Commissioner for Patents
 Box Patent Application
 Washington, DC 20231

5. ☐ Microfiche Computer Program (Appendix)
 6. Nucleotide and/or Amino Acid Sequence Submission (if applicable, all necessary).
 a. ☐ Computer Readable Copy
 b. ☐ Paper Copy (identical to computer copy)
 c. ☐ Statement verifying identity of above copies

ACCOMPANYING APPLICATION PARTS

7. ☒ Assignment Papers (cover sheet & document(s))
 8. ☐ 37 C.F.R. § 3.73(b) Statement of Power of Attorney (when there is an assignee)
 9. ☐ English Translation Document (if applicable)
 10. ☐ Information Disclosure Statement (IDS)/PTO-1449
 11. ☐ Preliminary Amendment
 12. ☒ Return Receipt Postcard (MPEP 503) (Should be specifically timed)
 13. ☐ * Small Entity Statement(s) Statement filed in prior application
 14. ☐ Status still proper and desired (PTO/SB/09-12)
 15. ☐ Certified Copy of Priority Document(s) (if foreign priority is claimed)
 16. ☐ Other:

16. If a CONTINUING APPLICATION, check appropriate box, and supply the requisite information below and in a preliminary amendment:
☐ Continuation ☐ Divisional ☐ Continuation-in-part (CIP) of prior application No. _____ / _____
 Prior application information: Examiner _____ Group / Art Unit _____

For CONTINUATION or DIVISIONAL APPS only: The entire disclosure of the prior application, from which an oath or declaration is supplied under Box 4b, is considered a part of the disclosure of the accompanying continuation or divisional application and is hereby incorporated by reference. The incorporation can only be relied upon when a portion has been inadvertently omitted from the submitted application parts.

17. CORRESPONDENCE ADDRESS

☐ Customer Number or Bar Code Label

(Insert Customer No. or Attach bar code label here)

or ☒ Correspondence address below

Name	Richard L. Byrne			
	Webb Ziesenheim Logsdon Orkin & Hanson, P.C.			
Address	700 Koppers Building			
	436 Seventh Avenue			
City	Pittsburgh	State	PA	Zip Code 15219
Country		Telephone	412-471-8815	Fax 412-471-4094

Name (Print/Type)	Richard L. Byrne	Registration No. (Attorney/Agent)	28,498
Signature	<i>Richard L. Byrne</i>	Date	11/15/00

Burden Hour Statement: This form is estimated to take 0.2 hours to complete. Time will vary depending upon the needs of the individual case. Any comments on the amount of time you are required to complete this form should be sent to the Chief Information Officer, Patent and Trademark Office, Washington, DC 20231. DO NOT SEND FEES OR COMPLETED FORMS TO THIS ADDRESS. SEND TO: Assistant Commissioner for Patents, Box Patent Application, Washington, DC 20231

FEE TRANSMITTAL

for FY 2000

Patent fees are subject to annual revision.

Small Entity payments must be supported by a small entity statement, otherwise large entity fees must be paid. See Forms PTO/SB/09-12. See 37 C.F.R. §§ 1.27 and 1.28.

TOTAL AMOUNT OF PAYMENT (\$1,092.00)

Complete if Known

Application Number
 Filing Date
 First Named Inventor Mark Kantrowitz
 Examiner Name
 Group / Art Unit
 Attorney Docket No. 2942-991842

METHOD OF PAYMENT (check one)

1. ☒ The Commissioner is hereby authorized to charge indicated fees and credit any overpayments to:

Deposit Account Number 23-0650

Deposit Account Name

- ☐ Charge Any Additional Fee Required Under 37 CFR §§ 1.16 and 1.17

2. ☒ Payment Enclosed:

☒ Check ☐ Money Order ☐ Other

FEE CALCULATION

1. BASIC FILING FEE

Large Entity Small Entity Fee Code (\$)	Fee Code (\$)	Fee Description	Fee Paid
101 690 201 345		Utility filing fee	710.00
106 310 206 155		Design filing fee	
107 480 207 240		Plant filing fee	
108 690 208 345		Reissue filing fee	
114 150 214 75		Provisional filing fee	

SUBTOTAL (1) (\$) 710.00

2. EXTRA CLAIM FEES

	Extra Claims	Fee from below	Fee Paid
Total Claims	39 - 20** = 19	X 18	= 342
Independent Claims	3 - 3** = 0	X 80	= 0
Multiple Dependent			= 0

**or number previously paid, if greater. For Reissues, see below

Large Entity Small Entity Fee Code (\$)	Fee Code (\$)	Fee Description
103 18 203 9		Claims in excess of 20
102 78 202 39		Independent claims in excess of 3
104 260 204 130		Multiple dependent claim, if not paid
109 78 209 39		** Reissue independent claims over original patent
110 18 210 9		** Reissue claims in excess of 20 and over original patent

SUBTOTAL (2) (\$) 342.00

FEE CALCULATION (continued)

3. ADDITIONAL FEES

Large Entity Small Entity Fee Code (\$)	Fee Code (\$)	Fee Description	Fee Paid
105 130 205 65		Surcharge - late filing fee or oath	0.00
127 50 227 25		Surcharge - late provisional filing fee or cover sheet	0.00
139 130 139 130		Non-English specification	0.00
147 2,520 147 2,520		For filing a request for reexamination	0.00
112 920* 112 920*		Requesting publication of SIR prior to Examiner action	0.00
113 1,840* 113 1,840*		Requesting publication of SIR after Examiner action	0.00
115 110 215 55		Extension for reply within first month	0.00
116 380 216 190		Extension for reply within second month	0.00
117 870 217 435		Extension for reply within third month	0.00
118 1,360 218 680		Extension for reply within fourth month	0.00
128 1,850 228 925		Extension for reply within fifth month	0.00
119 300 219 150		Notice of Appeal	0.00
120 300 220 150		Filing a brief in support of an appeal	0.00
121 260 221 130		Request for oral hearing	0.00
138 1,510 138 1,510		Petition to institute a public use proceeding	0.00
140 110 240 55		Petition to revive - unavoidable	0.00
141 1,210 241 605		Petition to revive - unintentional	0.00
142 1,210 242 605		Utility issue fee (or reissue)	0.00
143 430 243 215		Design issue fee	0.00
144 580 244 290		Plant issue fee	0.00
122 130 122 130		Petitions to the Commissioner	0.00
123 50 123 50		Petitions related to provisional applications	0.00
126 240 226 120		Submission of Information Disclosure Stmt	0.00
581 40 581 40		Recording each patent assignment per property (times number of properties)	40.00
146 690 246 345		Filing a submission after final rejection (37 CFR § 1.129(a))	0.00
149 690 249 345		For each additional invention to be examined (37 CFR § 1.129(b))	0.00
		Other fee (specify) _____	0.00
		Other fee (specify) _____	0.00

*Reduced by Basic Filing Fee Paid

SUBTOTAL (3) (\$) 40.00

SUBMITTED BY

Name (Print/Type) Richard L. Byrne

Registration No. (Attorney/Agent)

28,498

Complete (if applicable)

Telephone 471-8815

Signature *Richard L. Byrne*

Date 11/15/00

WARNING:

Information on this form may become public. Credit card information should not be included on this form. Provide credit card information and authorization on PTO-2038.

Burden Hour Statement: This form is estimated to take 0.2 hours to complete. Time will vary depending upon the needs of the individual case. Any comments on the amount of time you are required to complete this form should be sent to the Chief Information Officer, Patent and Trademark Office, Washington, DC 20231. DO NOT SEND FEES OR COMPLETED FORMS TO THIS ADDRESS. SEND TO: Assistant Commissioner for Patents, Washington, DC 20231.

PATENT APPLICATION

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

IN RE APPLICATION OF:

ATTORNEY'S DOCKET NUMBER

MARK KANTROWITZ

2942-991842

ENTITLED

"Method And Apparatus For Efficient Identification Of Duplicate And Near-Duplicate Documents And Text Spans Using High-Discriminability Text Fragments"

BOX PATENT APPLICATION

Assistant Commissioner for Patents
Washington, D.C. 20231

EXPRESS MAIL CERTIFICATE

"Express Mail" Label Number EL561512695US

Date of Deposit November 15, 2000

I hereby certify that the following attached papers or fee

UTILITY PATENT APPLICATION TRANSMITTAL (1 p.); FEE TRANSMITTAL FOR FY 2000 (1 p.); SPECIFICATION (12 pp.); CLAIMS (6 pp., 39 claims); ABSTRACT (1 p.); NINE SHEET OF DRAWINGS (Figs. 1-6); DECLARATION AND POWER OF ATTORNEY (2 pp.); RECORDATION FORM COVER SHEET - PATENTS ONLY (2 pp.); ASSIGNMENT (2 pp.); and two checks in amounts of \$1,052.00 and \$40.00

is being deposited with the United States Postal Service "Express Mail Post Office to Addressee" service under 37 C.F.R. §1.10 on the date indicated above and is addressed to the Assistant Commissioner for Patents, Washington, D.C. 20231.

Linda L. Marlowe

(Typed name of person mailing paper or fee)

Linda L. Marlowe

(Signature of person mailing paper or fee)

METHOD AND APPARATUS FOR EFFICIENT IDENTIFICATION OF
DUPLICATE AND NEAR-DUPLICATE DOCUMENTS AND TEXT SPANS
USING HIGH-DISCRIMINABILITY TEXT FRAGMENTS

BACKGROUND OF THE INVENTION

5 1. Field of the Invention

This invention relates to a computer-assisted method and apparatus for identifying duplicate and near-duplicate documents or text spans in a collection of documents or text spans, respectively.

10 2. Description of the Prior Art

The current art includes inventions that compare a single pair of known-to-be-similar documents to identify the differences between the documents. For example, the Unix "diff" program uses an efficient algorithm for finding the longest common sub-sequence (LCS) between two sequences, such as the lines in two documents. Aho, Hopcroft, and Ullman, *Data Structures and Algorithms*, Addison-Wesley Publishing Company, April 1987, pages 189-192. The lines that are left when the LCS is removed represent the changes needed to transform one document into another. Additionally, U.S. Patent No. 4,807,182 uses anchor points (points in common between two files) to identify differences between an original and a modified version of a document. There are also programs for comparing a pair of files, such as the Unix "cmp" program.

20 Another approach for comparing documents is to compute a checksum for each document. If two documents have the same checksum, they are likely to be identical. But comparing documents using checksums is an extremely fragile method, since even a single character change in a document yields a different checksum. Thus, checksums are good for identifying exact duplicates, but not for identifying near-duplicates. U.S. Patent No. 5,680,611 teaches the use of checksums to identify duplicate records. U.S. Patent No. 5,898,836 discloses the use of checksums to identify whether a region of a document has changed by comparing checksums for sub-document passages, for example, the text between HTML tags.

30 Patrick Juola's method, discussed in Juola, Patrick, *What Can We Do With Small Corpora? Document Categorization via Cross-Entropy*, Proceedings of Workshop on Similarity and Categorization, 1997, uses the average length of matching character n-grams (an n-gram is a string of characters that may comprise all or part of a word) to identify similar documents. For each window of consecutive characters in the

source document, the average length of the longest matching sub-sequence at each position in the target document is computed. This effectively computes the average length of match at each position within the target document (counting the number of consecutive matching characters starting from the first character of the *n*-gram) for every possible character *n*-gram within the source document. This technique depends on the frequency of the *n*-grams within the document by requiring the *n*-grams and all sub-parts (at least the prefix sub-parts) to be of high frequency. The Juola method focuses on applications involving very small training corpora, and has been applied to a variety of areas, including language identification, determining authorship of a document, and text classification. The method does not provide a measure of distinctiveness.

The prior art does not compare more than two documents, does not allow text fragments in each document to appear in a different or arbitrary order, is not selective in the choice of *n*-grams used to compare the documents, does not use the frequency of the *n*-grams across documents for selecting *n*-grams used to compare the documents, and does not permit a mixture of very low frequency and very high frequency components in the *n*-grams.

SUMMARY OF THE INVENTION

It is an object of the present invention to provide a method and apparatus for the efficient identification of duplicate and near-duplicate documents and text spans.

Accordingly, I have developed a method and apparatus for the efficient identification of duplicate and near-duplicate (i.e., substantially duplicate) documents and text spans which use high-discriminability text fragments for comparing documents.

Near-duplicate documents contain long stretches of identical text that are not present in other, non-duplicate documents. The long text fragments that are present in only a few documents (high-intermediate rarity) represent distinctive features that can be used to distinguish similar documents from dissimilar documents in a robust fashion. These text fragments represent a kind of "signature" for a document which can be used to match the document with near-duplicate documents and to distinguish the document from non-duplicate documents. Documents that overlap significantly on such text fragments will most likely be duplicates or near-duplicates. Overlap occurs not just when the text is excerpted, but also when deliberate changes have been made to the text, such as paraphrasing, interspersing comments by another author, and outright plagiarism.

Typically, as long as the document is not completely rewritten, there will be large text fragments that are specific to the document and its duplicates. On the other hand, text fragments in common between two non-duplicate documents will likely be in common with many other documents.

5 The present invention identifies duplicate and near-duplicate documents and text spans by identifying a small number of distinctive features for each document, for example, distinctive word n-grams likely to appear in duplicate or near-duplicate documents. The features act as a proxy for the full document, allowing the invention to compare documents by comparing their distinctive features. Documents having at least
10 one feature in common are compared with each other. Near-duplicate documents are identified by counting the proportion of the features in common between the two documents. Using these common features allows the present invention to find near-duplicate documents efficiently without needing to compare each document with all the other documents in the collection, for example, by pairwise comparison. By comparing
15 features instead of entire documents, the present invention is much faster in finding duplicate and near-duplicate documents in a large collection of documents than might be possible with prior document comparison algorithms.

 A key to the effectiveness of this method is the ability to find distinctive features. The features need to be rare enough to be common among only near-duplicate
20 documents, but not so rare as to be specific to just one document. An individual word may not be rare enough, but an n-gram containing the word might be. Longer n-grams might be too rare. Additionally, the distinctive features may include glue words (i.e., very common words) within the features but, preferably, not at either end. Thus, distinctive features may include words that are common to just a few documents and/or
25 words that are common to all but a few documents.

 Blindly gathering all n-grams of appropriate rarity would yield a computationally expensive algorithm. Thus, the number of distinctive features used must be small in order for the algorithm to be computationally efficient. The present invention incorporates several methods that strike a balance between appropriate rarity and
30 computational expense.

 Applications of the present invention include removing redundancy in document collections (including web catalogs and search engines), matching summary

sentences with corresponding document sentences, and detection of plagiarism and copyright infringement for text documents and passages.

BRIEF DESCRIPTION OF THE DRAWINGS

Fig. 1 is a flow diagram of a first embodiment of a method according to the present invention as applied to documents;

Fig. 2 is a flow diagram of a second embodiment of a method according to the present invention as applied to documents;

Figs. 3A and 3B are a flow diagram of a third embodiment of a method according to the present invention as applied to documents;

Fig. 3C is an illustration of a document index;

Fig. 3D is an illustration of a feature index;

Fig. 3E is an illustration of a list 324;

Fig. 3F is an illustration of a list 330;

Fig. 3G is an illustration of a list 336;

Fig. 4 is a flow diagram of an embodiment of a method according to the present invention as applied to text spans;

Fig. 5 is a flow diagram of an embodiment of a method according to the present invention as applied to images; and

Fig. 6 is an illustration of an apparatus according to the present invention.

DESCRIPTION OF THE PREFERRED EMBODIMENTS

Referring to Fig. 1, the present invention is utilized to find duplicate or near-duplicate documents within a document collection 100. Step 110 identifies distinctive features in the document collection 100 and in each document in the collection 100. Loop 112 iterates for each pair of documents. Within loop 112, step 114 determines if the pair of documents has at least one distinctive feature in common. If they do, the pair is compared in step 116 to determine if they are duplicate or near-duplicate documents. Loop 112 then continues with the next pair of documents. If the pair of documents does not have at least one distinctive feature in common, no comparison is performed, and loop 112 continues with the next pair of documents.

The method illustrated in Fig. 1 can be applied to, for example: removing duplicates in document collections; detecting plagiarism; detecting copyright infringement; determining the authorship of a document; clustering successive versions

of a document from among a collection of documents; seeding a text classification or text clustering algorithm with sets of duplicate or near-duplicate documents; matching an e-mail message with responses to the e-mail message, and vice versa; and creating a document index for use with a query system to efficiently find documents that contain a particular phrase or excerpt in response to a query, even if the particular phrase or excerpt was not recorded correctly in the document or the query.

The method can also be applied to augmenting information retrieval or text classification algorithms that use single-word terms with a small number of multi-word terms. Algorithms of this type that are based on a bag-of-words model assume that each word appears independently. Although such algorithms can be extended to apply to word bigrams, trigrams, and so on, allowing all word n-grams of a particular length rapidly becomes computationally unmanageable. The present invention may be used to generate a small list of word n-grams to augment the bag-of-words index. These word n-grams are likely to distinguish documents. Therefore, if they are present in a query, they can help narrow the search results considerably. This is in contrast to methods based on word co-occurrence statistics which yield word n-grams that are rather common in the document set.

The method illustrated in Fig. 1 may be used to determine whether documents are duplicates or near-duplicates even if the distinctive features appear in a different order in each document.

The distinctive features may be distinctive text fragments found within the collection of documents 100. As such, the method may be applied to information retrieval methods, such as a text classification method or any information retrieval method that assumes word independence and adds the distinctive text fragments to an index set.

The distinctive text fragments may be sequences of at least two words that appear in a limited number of documents in the document collection 100. If one distinctive text fragment is found within another distinctive text fragment, only the longest distinctive text fragment may be considered as a feature. A sequence of at least two words may be considered as appearing in a document when the document contains the sequence of at least two words at least a user-specified minimum number of times or

a user-specified minimum frequency. The frequency may be defined as the number of occurrences in the document divided by the length of the document.

5 For each sequence of at least two words, a distinctiveness score may be calculated and the highest scoring sequences that are found in at least two documents in the document collection 100 may be considered as text fragments. The distinctiveness score may be the reciprocal of the number of documents containing the phrase multiplied by a monotonic function of the number of words in the phrase, where the monotonic function may be the number of words in the phrase.

10 The limited number restricting the number of documents having the sequence of at least two words may be selected by a user as a constant or a percentage. The limited number may be defined by a linear function of the number of documents in the document collection 100, such as a linear function of the square root or logarithm of the number of documents in the document collection 100.

15 The distinctive text fragments may include glue words (i.e., words that appear in almost all of the documents and for which their absence is distinctive). Glue words include stopwords like "the" and "of" and allow phrases like "United States of America" to be counted as distinctive phrases. The method may exclude glue words that appear at either extreme of the distinctive text fragment. Again, the sequence of at least two words may be considered as appearing in a document when the document contains
20 the sequence of at least two words at least a user-specified minimum number of times or a user-specified minimum frequency. The frequency may be defined as the number of occurrences in the document divided by the length of the document.

Fig. 2 illustrates another embodiment of the present invention which finds duplicate or near-duplicate documents within a document collection 200. Step 210
25 identifies distinctive features of the documents in the document collection 200 and in each document in the collection 200. Loop 212 iterates for each pair of documents. Within loop 212, step 214 determines if the pair of documents has at least one distinctive feature in common. If they do, step 216 divides the number of features that the pair of documents has in common by the smaller number of the number of features in each document. Step 218 determines whether the result of step 216 is greater than a threshold value. The threshold value may be a constant, a fixed percentage of the number of documents in the document collection 200, the logarithm of the number of documents,
30

or the square root of the number of documents. If the result is greater than the threshold, step 220 deems the documents duplicates or near-duplicates, and loop 212 continues with the next pair of documents. If the result is not greater than the threshold, the documents are not duplicates or near-duplicates, and loop 212 continues.

5 Figs. 3A and 3B show another embodiment of the present invention which finds duplicate or near-duplicate documents within a document collection 300. Starting with Fig. 3A, step 310 identifies distinctive features of the documents in the document collection 300 and in each document in the collection 300. Step 312 builds a document index 314 and step 316 builds a feature index 318. The document index 314
10 maps each document to the features contained therein. The feature index 318 maps the features to the documents that contain them. The indexes 314 and 318 are built in a manner that ignores duplicates (i.e., if a feature is repeated within a document, it is mapped only once). Loop 320 iterates through each document such that step 322 can create a list 324 that includes each unique distinctive feature that was identified in step
15 310. For each distinctive feature in list 324, loop 326 iterates through the feature index 318 so that step 328 can create a list 330 that includes each distinctive feature and the documents in which the distinctive feature is located.

Referring now to Fig. 3B, loop 332 iterates through list 330. Within loop
332, step 334 creates a list 336 of pairs of documents that have at least one feature in
20 common and the number of features they have in common. Loop 338 iterates through list 336. For each pair of documents in list 336, step 340 divides the number of features that the pair of documents has in common by the smaller number of the number of features in each document (from the document index 314). Step 342 determines whether the result of step 340 is greater than a threshold value. The threshold value may, for
25 example, be a constant, a fixed percentage of the number of documents in the document collection 300, the logarithm of the number of documents, or the square root of the number of documents. If the result is greater than the threshold, step 344 deems the documents duplicates or near-duplicates, and loop 338 continues with the next pair of documents. If the result is not greater than the threshold, the documents are not
30 duplicates or near-duplicates, and loop 338 continues.

Fig. 3C illustrates an example format for the document index 314. Likewise, Fig. 3D illustrates the feature index 318, Fig. 3E illustrates list 324, Fig. 3F illustrates list 330, and Fig. 3G illustrates list 336 as constructed in two steps.

Referring to Fig. 4, a method according to the present invention is utilized to find duplicate or near-duplicate text spans, including sentences, within a text span collection 400. The text spans in the collection 400 may be sentences. Step 410 identifies distinctive features of the text spans in the text span collection 400 and in each text span in the collection 400. Loop 412 iterates for each pair of text spans. Within loop 412, step 414 determines if the pair of text spans has at least one distinctive feature in common. If they do, the pair is compared in step 416 to determine if they are duplicate or near-duplicate text spans. Loop 412 then continues with the next pair of text spans. If the pair of text spans does not have at least one distinctive feature in common, no comparison is performed, and loop 412 continues with the next pair of text spans.

This method may be used to match sentences from one document with sentences from another. This would be useful in matching sentences of a human-written summary for an original document with sentences from the original document. Similarly, in a plagiarism detector, once the method as applied to documents has found duplicate documents, the sentence version can be used to match sentences in the plagiarized copy with the corresponding sentences from the original document. Another application of sentence matching would identify changes made to a document in a word processing application where such changes need not retain the sentences, lines, or other text fragments in the original order.

Referring to Fig. 5, the present invention is utilized to find duplicate or near-duplicate images within an image collection 500. Step 510 identifies distinctive features of the images in the image collection 500 and in each image in the collection 500. The distinctive features may be sequences of at least two adjacent tiles from the images. Loop 512 iterates for each pair of images. Within loop 512, step 514 determines if the pair of images has at least one distinctive feature in common. If they do, the pair is compared in step 516 to determine if they are duplicate or near-duplicate images. Loop 512 then continues with the next pair of images. If the pair of images does not have at least one distinctive feature in common, no comparison is performed, and loop 512 continues with the next pair of images.

In a preferred embodiment of the invention according to the method illustrated in Fig. 5, the method performs canonicalization of the images by converting them to black and white and sampling them at several resolutions. As compared to the method applied to text, small overlapping tiles correspond to words and horizontal and vertical sequences to text fragments.

The method illustrated in Fig. 5 may be applied to detecting copyright infringement based on image content where the original image does not have a digital watermark. This method may also be applied to fingerprint identification or handwritten signature authentication, among other applications.

The present invention also includes an apparatus that is capable of identifying duplicate and near-duplicate documents in a large collection of documents. The apparatus includes a means for initially selecting distinctive features contained within the collection of documents, a means for subsequently identifying the distinctive features contained in each document, and a means for then comparing the distinctive features of each pair of documents having at least one distinctive feature in common to determine whether the documents are duplicate or near-duplicate documents.

Fig. 6 illustrates an embodiment of an apparatus of the present invention capable of enabling the methods of the present invention. A computer system 600 is utilized to enable the method. The computer system 600 includes a display unit 610 and an input device 612. The input device 612 may be any device capable of receiving user input, for example, a keyboard or a scanner. The computer system 600 also includes a storage device 614 for storing the document collection and a storage device 616 for storing the method according to the present invention. A processor 618 executes the method stored on storage device 616 and accesses the document collection stored on storage device 614. The processor is also capable of sending information to the display unit 610 and receiving information from the input device 612. Any type of computer system having a variety of software and hardware components which is capable of enabling the methods according to the present invention may be used, including, but not limited to, a desktop system, a laptop system, or any network system.

The present invention was implemented in accordance with the method illustrated in Figs. 3A and 3B. In the implementation, $DF(x)$ was the number of documents containing the text "x", N was the overall number of documents, and R was

a threshold on DF. Possible choices for R included a constant, a fixed percentage of N (for example five percent), the logarithm of N, or the square root of N.

5 A first pass over all the documents computed $DF(x)$ for all words in the documents after converting the words to lowercase and removing punctuation from the beginning and end of the word. Optionally, a word in a particular document may be restricted from contributing to $DF(x)$ if the word's frequency in that document falls below a user-specified threshold.

10 A second pass gathered the distinctive features or phrases. A phrase consisted of at least two words which occur in more than one document and in no more than R documents ($1 < DF(x) < R$). The phrases also contained glue words that occurred in at least $(N - R)$ documents. The glue words could appear within a phrase, but not in the leftmost or rightmost position in the phrase. Essentially, the document was segmented at words of intermediate rarity ($R < DF(x) < R - N$) and what remained were considered distinctive phrases. Optionally, the phrases may also be segmented at the glue
15 words to obtain additional distinctive sub-phrases, for example, "United States of America" yields "United States" upon splitting at the "of". The second pass also built a document index that mapped each document to its set of distinctive phrases and sub-phrases using a document identifier and a phrase identifier and built a phrase index that mapped from the phrases to the documents that contained them using the phrase and
20 document identifiers. The indexes were built in a manner that ignores duplicates.

Unlike single words of low DF, the phrases were long enough to distinguish documents that happened to use the same vocabulary, but short enough to be common among duplicate documents.

25 A third pass iterated over the document identifiers in the document index (it is not necessary to use the actual documents once the indexes are built). For each document identifier, the document index was used to gather a list of the phrase identifiers. For each phrase identifier, the document identifiers obtained from the phrase index was iterated over to count the total number of times each document identifier occurred. Thus, for each document identifier, a list of documents that overlap with the
30 document in at least one phrase and the number of phrases that overlap was generated. This list of document identifiers included only those documents that had at least one phrase in common with the source document in order to avoid the need to compare the

source document with every other document. For each pair of documents, an overlap ratio was calculated by dividing the number of common phrases by the smaller of the number of phrases in each document. This made it possible to detect a small passage excerpted from a longer document. The overlap ratio was compared with a match percentage threshold. If it exceeded the threshold, the pair was reported as potential near-duplicates. Optionally, the results may be accepted as is or a more detailed comparison algorithm may be applied to the near-duplicate document pairs.

This implementation is rather robust since small changes to a document have little impact on the effectiveness of the method. If there are any telltale signs of the original document left, this method will find them. Moreover, the distinctive phrases do not need to appear in the same order in the duplicate documents.

The implementation is also very efficient. The first two passes are linear in N . The third pass runs in time $N * P$, where P is the average number of documents that overlap in at least one phrase. In the worst case P is N , but typically P is R . Note that as R increases, so does the accuracy, but the running time also increases. So, there is a trade-off between running time and accuracy. In practice, however, an acceptable level of accuracy is achieved for a running time that is linear in N . This is a significant improvement over algorithms which would require pairwise comparisons of all the documents, or at least N -squared running time.

The implementation was executed on 125 newspaper articles and their corresponding human-written summaries, for a total of 250 documents. For each pair of documents identified as near-duplicates, if the pair consisted of an article and its summary, it was counted as a correct match. Otherwise, it was counted as an incorrect match. For the purpose of the experiment, pairs consisting of a document and itself were excluded because the implementation successfully matches any document with itself. Using a minimum overlap threshold of 25% and a DF threshold of 5%, the method processed all 250 documents in 13 seconds and was able to match 232 of the 250 documents with their corresponding summary or article correctly, and none incorrectly. This represents a precision (accuracy) of 100%, a recall (coverage) of 92.8%, and an F1 score (harmonic mean of precision and recall) of 96.3%. Inspection of the results showed that in all cases where the algorithm did not find a match, the highest ranking document, although below the threshold, was the correct match.

The implementation may use different thresholds for the low frequency and glue words. Sequences of mid-range DF words where the sequence itself has low DF, may be included. Additionally, the number of words in a phrase may be factored in as a measure of the phrase's complexity in addition to rarity, for example, dividing the length of the phrase by the phrase's DF (TL/DF or $\log(TL)/DF$). Although, this yields a preference for longer phrases, it allows longer phrases to have higher DF and, thus, be less distinctive.

It will be understood by those skilled in the art that while the foregoing description sets forth in detail preferred embodiments of the present invention, modifications, additions, and changes may be made thereto without departing from the spirit and scope of the invention. Having thus described my invention with the detail and particularity required by the Patent Laws, what is desired to be protected by Letters Patent is set forth in the following claims.

1 claim:

1. A computer-assisted method for identifying duplicate and near-duplicate documents in a large collection of documents, comprising the steps of:
 - initially, selecting distinctive features contained in the collection of documents,
 - 5 then, for each document, identifying the distinctive features contained in the document, and
 - then, for each pair of documents having at least one distinctive feature in common, comparing the distinctive features of the documents to determine whether the documents are duplicate or near-duplicate documents.
2. The computer-assisted method according to claim 1, wherein the method is applied to removing duplicates in document collections.
3. The computer-assisted method according to claim 1, wherein the method is applied to detecting plagiarism.
4. The computer-assisted method according to claim 1, wherein the method is applied to detecting copyright infringement.
5. The computer-assisted method according to claim 1, wherein the method is applied to determine the authorship of a document.
6. The computer-assisted method according to claim 1, wherein the method is applied to clustering successive versions of a document from among a collection of documents.
7. The computer-assisted method according to claim 1, wherein the method is applied to seeding a text classification or text clustering algorithm with sets of duplicate or near-duplicate documents.

8. The computer-assisted method according to claim 1, wherein the method is applied to matching an e-mail message with responses to the e-mail message.

9. The computer-assisted method according to claim 1, wherein the method is applied to matching responses to an e-mail message with the e-mail message.

10. The computer-assisted method according to claim 1, wherein the method is applied to creating a document index for use with a query system to efficiently find documents in response to a query which contain a particular phrase or excerpt.

11. The computer-assisted method according to claim 10, wherein the document index can be utilized even if the particular phrase or excerpt was not recorded correctly in the document or in the query.

12. The computer-assisted method according to claim 1, wherein the distinctive features appear in a different order in each of the documents.

13. The computer-assisted method according to claim 1, wherein the distinctive features are distinctive text fragments from the documents in the document collection.

14. The computer-assisted method according to claim 13, wherein the method is applied to information retrieval methods.

15. The computer-assisted method according to claim 14, wherein the information retrieval method is a text classification method.

16. The computer-assisted method according to claim 14, wherein: the information retrieval method assumes word independence, and the distinctive text fragments are added to an index set.

17. The computer-assisted method according to claim 13, wherein the distinctive text fragments are sequences of at least two words that appear in a limited number of documents in the document collection.

18. The computer-assisted method according to claim 14, wherein if one distinctive text fragment is contained within another distinctive text fragment within the same document, only the longest distinctive text fragment is considered as a distinctive feature.

19. The computer-assisted method according to claim 17, wherein the sequences of at least two words are considered as appearing in a document when the document contains the sequence of at least two words at least a user-specified minimum number of times.

20. The computer-assisted method according to claim 17, wherein the sequences of at least two words are considered as appearing in a document when the document contains the sequence of at least two words at least a user-specified minimum frequency.

21. The computer-assisted method according to claim 17, wherein:
for each sequence of at least two words, a distinctiveness score is calculated, and

the highest scoring sequences that are found in at least two documents in
5 the document collection are considered distinctive text fragments.

22. The computer-assisted method according to claim 21, wherein the distinctiveness score is the reciprocal of the number of documents containing the phrase multiplied by a monotonic function of the number of words in the phrase.

23. The computer-assisted method according to claim 21, wherein the monotonic function is the number of words in the phrase.

24. The computer-assisted method according to claim 21, wherein the distinctiveness score is the percentage of documents not containing the phrase multiplied by a monotonic function of the number of words in the phrase.

25. The computer-assisted method according to claim 24, wherein the monotonic function is the number of words in the phrase.

26. The computer-assisted method according to claim 17, wherein the limited number is selected by a user.

27. The computer-assisted method according to claim 17, wherein the limited number is defined by a linear function of the number of documents in the document collection.

28. The computer-assisted method according to claim 17, wherein the distinctive text fragments include glue words.

29. The computer-assisted method according to claim 28, wherein the glue words do not appear at either extreme of the distinctive text fragments.

30. The computer-assisted method according to claim 1, further including the step of for each pair of documents having at least one distinctive feature in common, counting the number of distinctive features in common,

wherein determining whether the pair of documents is duplicates or near-
5 duplicates includes the steps of:

for each pair of documents, calculating an overlap ratio by dividing the number of distinctive features in common by the smaller of the number of distinctive features per document, and

10 comparing the overlap ratio to a threshold and if the overlap ratio is greater than the threshold, then the pair of documents are duplicates or near-duplicates, otherwise the pair of documents is not duplicates or near-duplicates.

31. The computer-assisted method according to claim 30, further including the steps of:

building a document index that maps each document to its associated distinctive features, wherein if one distinctive feature is repeated within one document,
5 the index maps the document to the distinctive feature once, and

building a feature index that maps each distinctive feature to its associated document, wherein if one distinctive feature is repeated within one document, the index maps the distinctive feature to the document once,

wherein determining whether the pair of documents are duplicates or
10 near-duplicates further includes the steps of:

creating a list of unique distinctive features from the document index,

for each unique distinctive feature, creating a list of documents which contain the unique distinctive feature, and

15 for each document, creating a list of documents that have at least one feature in common with the document and the number of features in common with the document.

32. The computer-assisted method according to claim 31, wherein the distinctive features include distinctive phrases.

33. The computer-assisted method according to claim 31, wherein the distinctive features appear in a different order in each of the documents.

34. The computer-assisted method according to claim 31, wherein the distinctive features include text spans.

35. The computer-assisted method according to claim 34, wherein the text spans include sentences.

36. The computer-assisted method according to claim 34, wherein the text spans include lines of text.

37. A computer-assisted method for identifying duplicate and near-duplicate text spans in a large collection of text spans, comprising the steps of:

initially, selecting distinctive features contained in the collection of text spans,

5 then, for each text span, identifying the distinctive features contained in the text span, and

then, for each pair of text spans having at least one distinctive feature in common, comparing the distinctive features of the text spans to determine whether the text spans are duplicate or near-duplicate text spans.

38. The computer-assisted method according to claim 37, wherein the text spans are sentences.

39. An apparatus to enable a method for identifying duplicate and near-duplicate documents in a large collection of documents, comprising:

a means for initially selecting distinctive features contained in the collection of documents;

5 a means for subsequently identifying the distinctive features contained in each document; and

a means for then comparing the distinctive features of each pair of documents having at least one distinctive feature in common to determine whether the documents are duplicate or near-duplicate documents.

METHOD AND APPARATUS FOR EFFICIENT IDENTIFICATION OF
DUPLICATE AND NEAR-DUPLICATE DOCUMENTS AND TEXT SPANS
USING HIGH-DISCRIMINABILITY TEXT FRAGMENTS

ABSTRACT OF THE DISCLOSURE

5 Disclosed is a computer-assisted method for finding duplicate or near-
duplicate documents or text spans within a document collection by using high-
discriminability text fragments. Distinctive features of the documents or text spans are
identified. For each pair of documents or text spans with at least one distinctive feature
in common, the distinctive features of each document or text span are compared to
10 determine whether the pair is duplicates or near-duplicates. An apparatus for performing
this computer-assisted method is also disclosed.

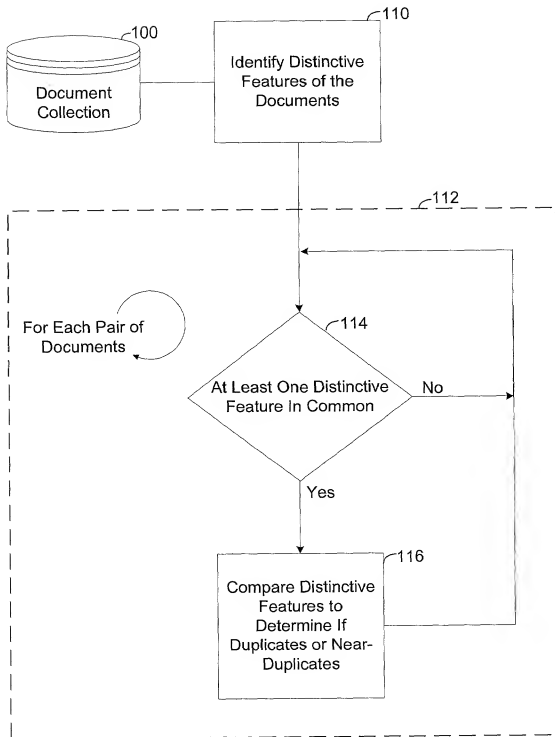


FIG. 1

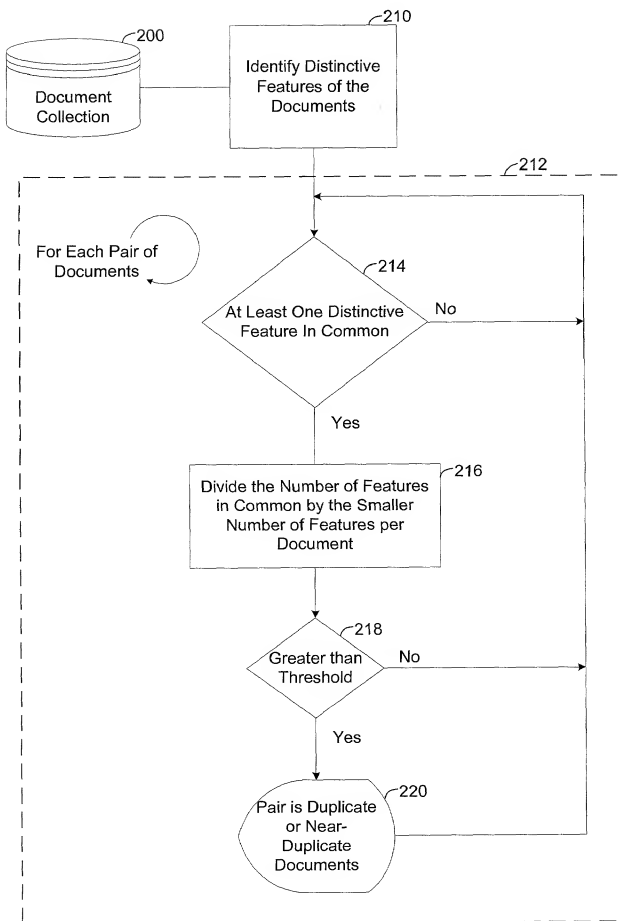


FIG. 2

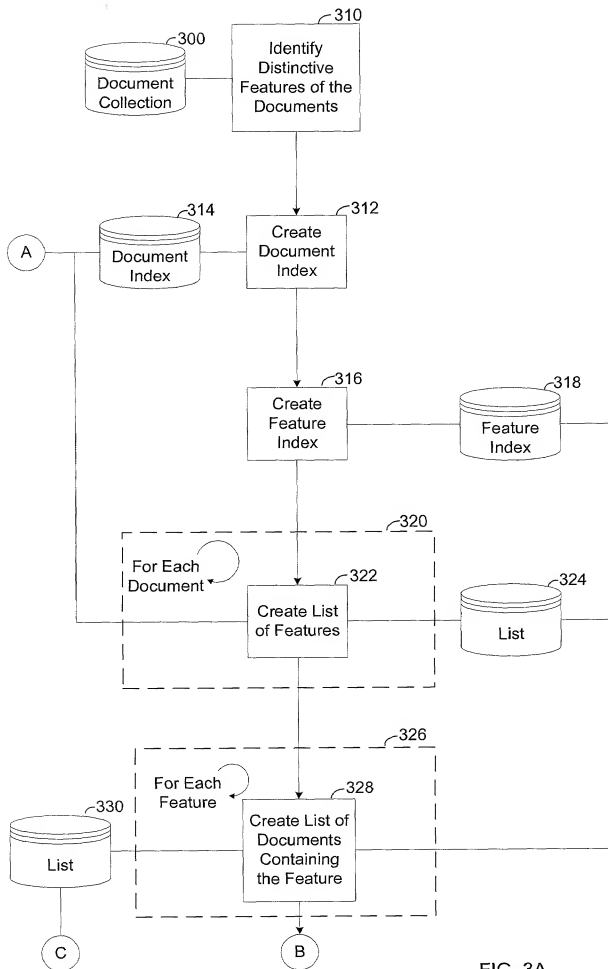


FIG. 3A

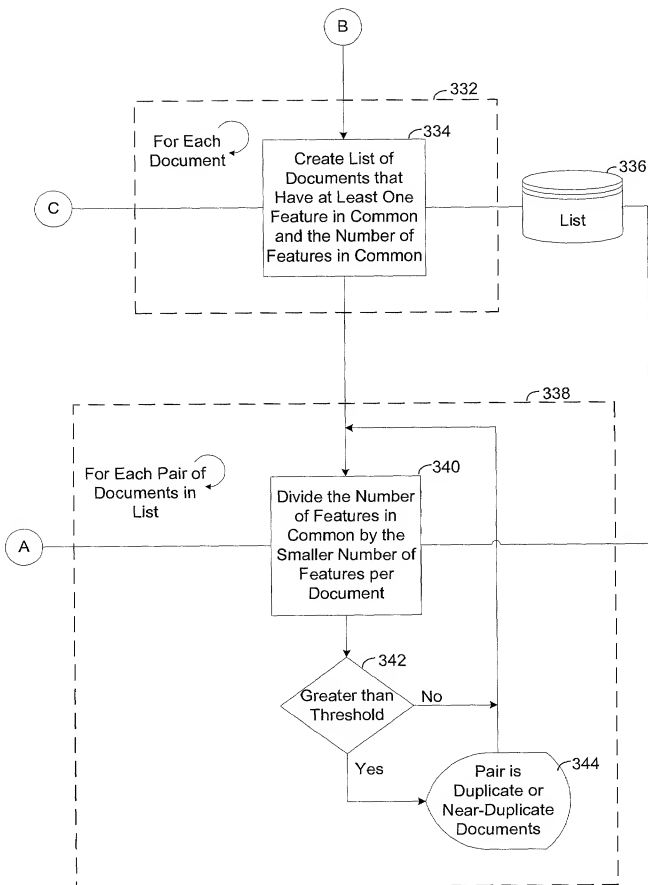


FIG. 3B

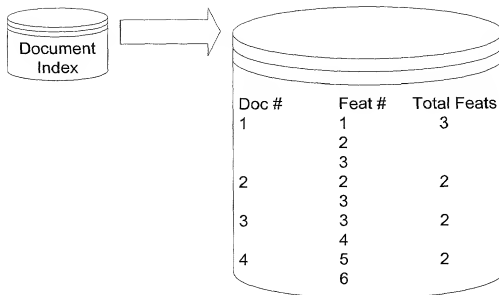


FIG. 3C

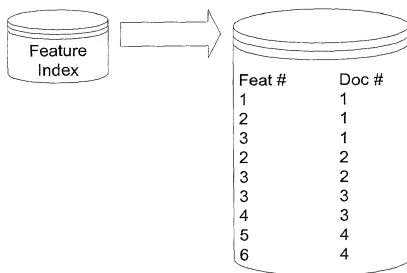


FIG. 3D

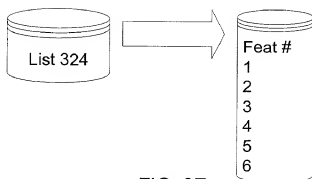


FIG. 3E

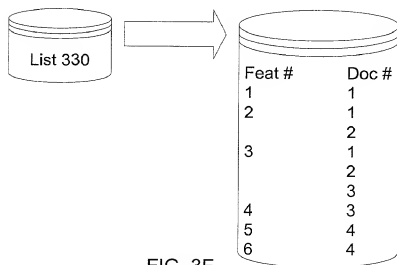


FIG. 3F

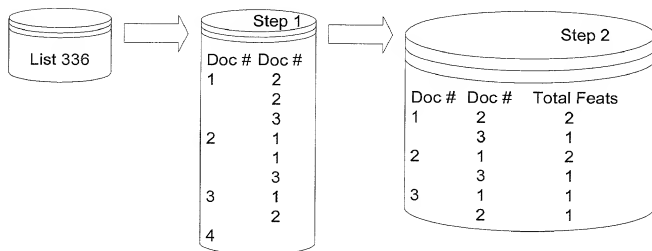


FIG. 3G

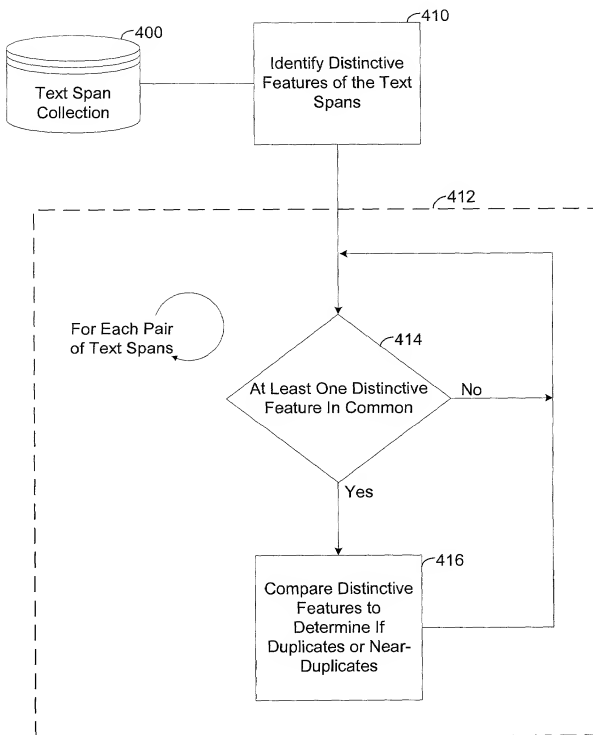


FIG. 4

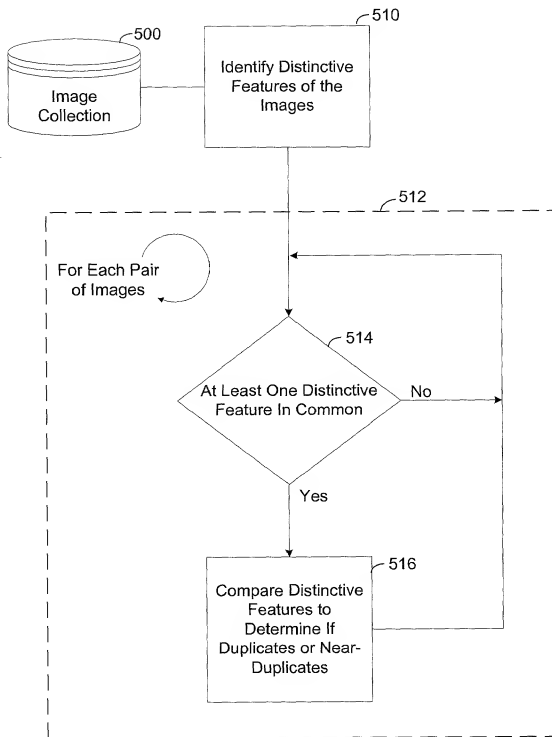


FIG. 5

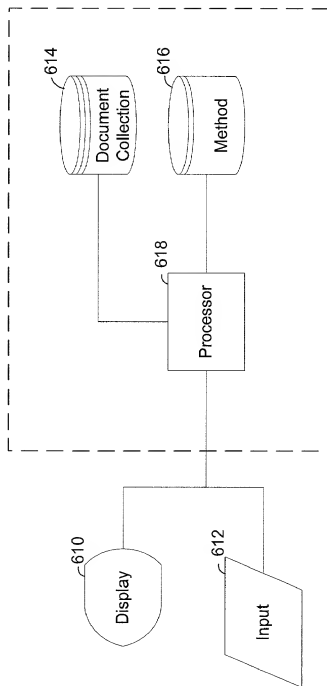


FIG. 6

DECLARATION AND POWER OF ATTORNEY

MARK KANTROWITZ, declares:

I am a citizen of the United States of America and a resident of **PITTSBURGH, ALLEGHENY COUNTY, COMMONWEALTH OF PENNSYLVANIA**, whose post-office address is **1363 SHADY AVENUE, PITTSBURGH, PA 15217**.

I believe myself to be the original, first and sole inventor of the improvement entitled **METHOD AND APPARATUS FOR EFFICIENT IDENTIFICATION OF DUPLICATE AND NEAR-DUPLICATE DOCUMENTS AND TEXT SPANS USING HIGH-DISCRIMINABILITY TEXT FRAGMENTS** which is described and claimed in the annexed specification.

I have reviewed and understand the contents of the specification, including the claims.

I do not know and do not believe that the same was ever known or used in the United States before my invention thereof; or patented or described in any printed publication in any country before my invention or more than one year prior to this application; or in public use or on sale in the United States more than one year prior to this application.

Said invention has not been patented or been made the subject of an inventor's certificate in any country foreign to the United States on an application filed by me or my legal representatives or assigns more than twelve months prior to this application.

I acknowledge my duty to disclose information of which I am aware which is material to the patentability of this application in accordance with Title 37, Code of Federal Regulations, § 1.56(a).

No application for patent or inventor's certificate thereon has been filed by me or my legal representatives or assigns in any country foreign to the United States.

I declare further that all statements made herein of my own knowledge are true and that all statements made on information and belief are believed to be true; and further that these statements were made with the knowledge that willful false statements and the like so made are punishable by fine or imprisonment or both, under Section 1001 of Title 18 of the United States Code and that such willful false statements may jeopardize the validity of the application or any patent issuing thereon.

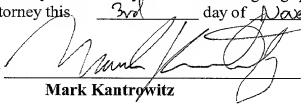
I hereby appoint William H. Logsdon, Registration No. 22,132; Russell D. Orkin, Registration No. 25,363; David C. Hanson, Registration No. 23,024; Richard L. Byrne, Registration No. 28,498; Frederick B. Ziesenheim, Registration No. 19,438; Kent E. Baldauf, Registration No. 25,826; Barbara E. Johnson, Registration No. 31,198; Paul M. Reznick, Registration No. 33,059; John W. McIlvaine, Registration No. 34,219; Michael I. Shamos, Registration No. 30,424; Blynn L. Shideler, Registration No. 35,034; Julie W. Meder, Registration No. 36,216; Lester N. Fortney, Registration No. 38,141; Randall A. Notzen, Registration No. 36,882; James G. Porcelli, Registration No. 33,757; Kent E. Baldauf, Jr., Registration No. 36,082; Christian E. Schuster, Registration No. 43,908; ; Dean E. Geibel,

Registration No. 42,570; Thomas J. Clinton, Registration No. 40,561; Nathan J. Prepelka, Registration No. 43,016; Jessica M. Sosenko, Registration No. P-47,102; Kirk M. Miles, Registration No. 37,891; and J. Matthew Pritchard, Registration No. 46,228, whose post-office address is 700 Koppers Building, 436 Seventh Avenue, Pittsburgh, Pennsylvania 15219-1818, Telephone No. 412-471-8815, my attorneys with full power of substitution and revocation, to prosecute this application and to transact all business in the Patent and Trademark Office connected therewith, to amend the specification, to appeal in case of rejection, as they may deem advisable, to receive the patent when granted and generally to do all matters and things needful in the premises, as fully and to all intents and purposes as I could do.

All correspondence and telephone calls should be addressed to Richard L. Byrne.

I hereby subscribe my name to the foregoing specification and claims, declaration and power of attorney this 3rd day of November, 2000.

Inventor


Mark Kantrowitz